

US009754569B2

(12) **United States Patent**  
**Neuhauser et al.**

(10) **Patent No.:** **US 9,754,569 B2**  
(45) **Date of Patent:** **\*Sep. 5, 2017**

(54) **AUDIO MATCHING WITH SEMANTIC AUDIO RECOGNITION AND REPORT GENERATION**

(71) Applicant: **The Nielsen Company (US), LLC**,  
New York, NY (US)

(72) Inventors: **Alan Neuhauser**, Silver Springs, MD  
(US); **John Stavropoulos**, Edison, NJ  
(US)

(73) Assignee: **THE NIELSEN COMPANY (US), LLC**,  
New York, NY (US)

(\*) Notice: Subject to any disclaimer, the term of this  
patent is extended or adjusted under 35  
U.S.C. 154(b) by 27 days.  
  
This patent is subject to a terminal dis-  
claimer.

(21) Appl. No.: **14/877,296**

(22) Filed: **Oct. 7, 2015**

(65) **Prior Publication Data**

US 2016/0027418 A1 Jan. 28, 2016

**Related U.S. Application Data**

(63) Continuation of application No. 13/725,004, filed on  
Dec. 21, 2012, now Pat. No. 9,183,849, which is a  
(Continued)

(51) **Int. Cl.**  
**G10L 15/00** (2013.01)  
**G10H 1/00** (2006.01)  
(Continued)

(52) **U.S. Cl.**  
CPC ..... **G10H 1/0008** (2013.01); **G10L 25/81**  
(2013.01); **G10L 25/84** (2013.01);  
(Continued)

(58) **Field of Classification Search**  
CPC ..... G10L 15/04; G10L 25/03; G10L 13/10;  
G06F 17/28  
(Continued)

(56) **References Cited**

U.S. PATENT DOCUMENTS

2,662,168 A 12/1953 Scherbatskoy  
3,919,479 A 11/1975 Moon et al.  
(Continued)

FOREIGN PATENT DOCUMENTS

EP 0887958 12/1998  
WO 9111062 7/1991  
(Continued)

OTHER PUBLICATIONS

IP Australia, "Notice of Acceptance," issued in connection with  
Australian Patent Application No. 2013361099, Apr. 21, 2016 (2  
pages).

(Continued)

*Primary Examiner* — Daniel Abebe

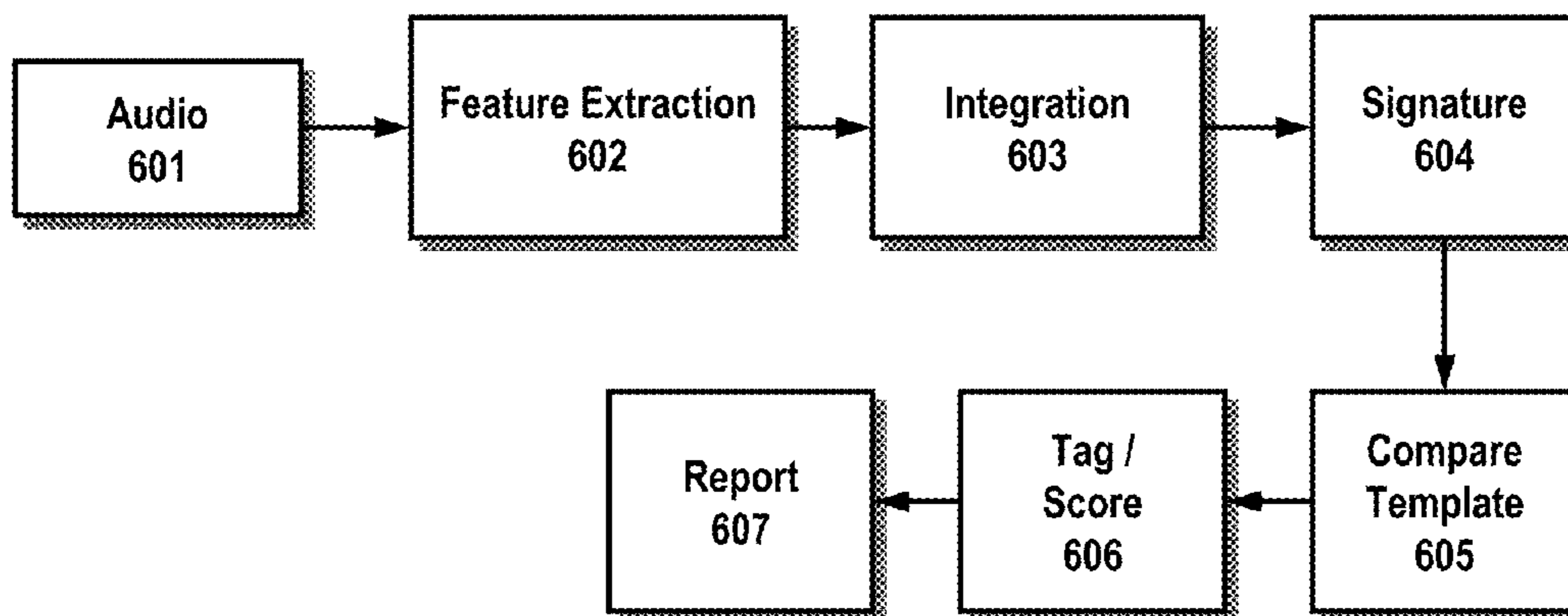
(74) *Attorney, Agent, or Firm* — Hanley, Flight &  
Zimmerman, LLC

(57) **ABSTRACT**

System, apparatus and method for determining semantic  
information from audio, where incoming audio is sampled  
and processed to extract audio features, including temporal,  
spectral, harmonic and rhythmic features. The extracted  
audio features are compared to stored audio templates that  
include ranges and/or values for certain features and are  
tagged for specific ranges and/or values. The semantic  
information may be associated with audio signature data

Extracted audio features that are most similar to one or more  
templates from the comparison are identified according to  
the tagged information. The tags are used to determine the  
semantic audio data that includes genre, instrumentation,  
style, acoustical dynamics, and emotive descriptor for the  
audio signal.

**20 Claims, 9 Drawing Sheets**



**Related U.S. Application Data**

continuation-in-part of application No. 13/724,836,  
filed on Dec. 21, 2012, now Pat. No. 9,195,649.

(51) **Int. Cl.**

*G10L 25/81* (2013.01)

*G10L 25/84* (2013.01)

*G10L 25/51* (2013.01)

(52) **U.S. Cl.**

CPC . *G10H 2210/036* (2013.01); *G10H 2210/041*  
(2013.01); *G10H 2210/056* (2013.01); *G10H*  
*2210/066* (2013.01); *G10H 2210/071*  
(2013.01); *G10H 2210/076* (2013.01); *G10H*  
*2210/081* (2013.01); *G10H 2250/235*  
(2013.01); *G10H 2250/251* (2013.01); *G10L*  
*15/00* (2013.01); *G10L 25/51* (2013.01)

(58) **Field of Classification Search**

USPC ..... 704/203  
See application file for complete search history.

(56)

**References Cited**

U.S. PATENT DOCUMENTS

4,230,990	A	10/1980	Lert, Jr. et al.
4,450,531	A	5/1984	Kenyon et al.
4,677,466	A	6/1987	Lert, Jr. et al.
4,697,209	A	9/1987	Kiewit et al.
4,739,398	A	4/1988	Thomas et al.
4,843,562	A	6/1989	Kenyon et al.
4,918,730	A	4/1990	Schulze
4,955,070	A	9/1990	Welsh et al.
5,436,653	A	7/1995	Ellis et al.
5,450,490	A	9/1995	Jensen et al.
5,457,768	A	10/1995	Tsuboi et al.
5,512,933	A	4/1996	Wheatley et al.
5,574,962	A	11/1996	Fardeau et al.
5,579,124	A	11/1996	Aijala et al.
5,581,800	A	12/1996	Fardeau et al.
5,594,934	A	1/1997	Lu et al.
5,612,729	A	3/1997	Ellis et al.
5,764,763	A	6/1998	Jensen et al.
5,787,334	A	7/1998	Fardeau et al.
5,918,223	A	6/1999	Blum et al.
6,201,176	B1	3/2001	Yourlo
6,574,594	B2	6/2003	Pitman et al.
6,604,072	B2	8/2003	Pitman et al.
6,675,174	B1	1/2004	Bolle et al.
6,871,180	B1	3/2005	Neuhauser et al.
6,968,564	B1	11/2005	Srinivasan
6,973,574	B2	12/2005	Mihcak et al.
6,990,453	B2	1/2006	Wang et al.
7,003,515	B1	2/2006	Glaser et al.
7,031,921	B2	4/2006	Pitman et al.
7,091,409	B2	8/2006	Li et al.
7,174,293	B2	2/2007	Kenyon et al.
7,284,255	B1	10/2007	Apel et al.
7,532,943	B2	5/2009	Weare
7,582,823	B2	9/2009	Kim et al.
7,634,406	B2	12/2009	Li et al.
7,640,141	B2	12/2009	Kolessar et al.
7,647,604	B2	1/2010	Ramaswamy
7,783,489	B2	8/2010	Kenyon et al.
7,982,117	B2	7/2011	Alcalde et al.
8,140,331	B2	3/2012	Lou
8,769,294	B2	7/2014	Samari et al.
8,825,188	B2	9/2014	Stone et al.
8,892,565	B2	11/2014	Xu et al.
8,959,016	B2	2/2015	McKenna et al.
9,640,156	B2	5/2017	Neuhauser et al.
2002/0181711	A1	12/2002	Logan et al.
2005/0177361	A1	8/2005	Srinivasan
2005/0232411	A1	10/2005	Srinivasan et al.
2007/0250777	A1	10/2007	Chen et al.

2007/0276667	A1	11/2007	Atkin et al.
2008/0032622	A1	2/2008	Kopra et al.
2008/0195654	A1	8/2008	Weare
2009/0277322	A1	11/2009	Cai et al.
2009/0306797	A1	12/2009	Cox et al.
2009/0313019	A1	12/2009	Kato et al.
2010/0161315	A1	6/2010	Melamed et al.
2010/0212478	A1	8/2010	Taub et al.
2011/0075851	A1	3/2011	LeBoeuf et al.
2011/0161076	A1	6/2011	Davis et al.
2012/0203363	A1	8/2012	McKenna et al.
2014/0019138	A1	1/2014	Stephens, Jr.
2014/0056432	A1	2/2014	Loui et al.
2014/0056433	A1	2/2014	Emerson, III
2014/0180673	A1	6/2014	Neuhauser et al.
2014/0180675	A1	6/2014	Neuhauser et al.
2014/0195221	A1	7/2014	Frank et al.
2014/0376729	A1*	12/2014	Crockett ..... G10L 15/04 381/56
2015/0332669	A1*	11/2015	Kalampoukas ..... G10L 25/27 704/243
2016/0012807	A1	1/2016	Neuhauser et al.
2016/0035332	A1	2/2016	Neuhauser et al.

FOREIGN PATENT DOCUMENTS

WO	0211123	2/2002
WO	03091990	11/2003
WO	2012168740	12/2012

OTHER PUBLICATIONS

Canadian Intellectual Property Office, "Examination Report," issued in connection with Canadian Patent Application No. 2,896,096, mailed May 9, 2016 (5 pages).

IP Australia, "Patent Examination Report No. 1," issued in connection with Application No. 2013361099, Feb. 26, 2016, 2 pages.

Patent Cooperation Treaty, "International Search Report," issued by the International Searching Authority in connection with PCT application No. PCT/US2013/076934, mailed Apr. 22, 2014 (6 pages).

Qing et al., "A Probabilistic Music Recommender Considering User Opinions and Audio Features," Information Processing and Management, 2007, pp. 473-487, vol. 43, (15 pages).

Tsekeridou et al., "Content-Based Video Parsing and Indexing Based on Audio-Visual Interaction," IEEE Transactions on Circuits and Systems for Video Technology, Apr. 2001, vol. 11, No. 4 (14 pages).

Haitsma et al., "A Highly Robust Audio Fingerprinting System," 2002, Philips Research, (9 pages).

Wold et al., "Content-Based Classification, Search, and Retrieval of Audio," Muscle Fish, 1996, IEEE (10 pages).

Belletini, et al., "A Framework for Robust Audio Fingerprinting," Journal of Communications, vol. 5, No. 5, Academy Publisher, May 2010 (16 pages).

Cano et al., "A Review of Algorithms for Audio Fingerprinting," 2002, IEEE (5 pages).

United States Patent and Trademark Office, "Office Action," issued in connection with U.S. Appl. No. 13/724,836 on Oct. 29, 2014 (6 pages).

Harte et al., "Detecting Harmonic Change in Musical Audio," AMCMM '06 Proceedings of the 1st ACM workshop on Audio and music computing multimedia, 2006, pp. 21-26, New York, NY, USA, (5 pages).

Klapuri, A., "Sound Onset Detection by Applying Psychoacoustic Knowledge," Acoustics, Speech, and Signal Processing, Proceedings, IEEE International Conference, 1999, pp. 3089-3092, vol. 6. IEEE, Phoenix, AZ, (4 pages).

Wang, Avery Li-Chun, "An Industrial-Strength Audio Search Algorithm," Proceedings of the International Conference on Music Information Retrieval (ISMIR), 2003, pp. 7-13, Baltimore, USA, (7 pages).



(56)

**References Cited**

OTHER PUBLICATIONS

United States Patent and Trademark Office, "Notice of Allowance," issued in connection with U.S. Appl. No. 13/725,021, on Jun. 5, 2015 (29 pages).

United States Patent and Trademark Office, "Notice of Allowance," issued in connection with U.S. Appl. No. 13/724, 836, on Mar. 30, 2015 (13 pages).

United States Patent and Trademark Office, "Supplemental Notice of Allowability," issued in connection with U.S. Appl. No. 13/724, 836, on May 14, 2015 (2 pages).

United States Patent and Trademark Office, "Notice of Allowance," issued in connection with U.S. Appl. No. 13/724, 836, on Jul. 21, 2015 (14 pages).

Patent Cooperation Treaty, "International Preliminary Report on Patentability," issued by the International Searching Authority in connection with PCT application No. PCT/US2013/076934, Jun. 23, 2015 (9 pages).

United States Patent and Trademark Office, "Corrected Notice of Allowability," issued in connection with U.S. Appl. No. 13/724, 836, Sep. 14, 2015, 6 pages.

United States Patent and Trademark Office, "Notice of Allowability," issued in connection with U.S. Appl. No. 13/725,021, Sep. 15, 2015, 10 pages.

United States Patent and Trademark Office, "Supplemental Notice of Allowability," issued in connection with U.S. Appl. No. 13/724,836, Sep. 30, 2015, 6 pages.

United States Patent and Trademark Office, "Notice of Allowance", issued in connection with U.S. Appl. No. 13/725,004, dated Oct. 7, 2015 (2 pages).

United States Patent and Trademark Office, "Notice of Allowance", issued in connection with U.S. Appl. No. 13/725,004, dated Sep. 17, 2015 (5 pages).

United States Patent and Trademark Office, "Notice of Allowance", issued in connection with U.S. Appl. No. 13/725,004, dated May 29, 2015 (7 pages).

United States Patent and Trademark Office, "Non-Final Office Action", issued in connection with U.S. Appl. No. 13/725,004, dated Nov. 19, 2014 (9 pages).

United States Patent and Trademark Office, "Supplemental Notice of Allowability," issued in connection with U.S. Appl. No. 13/725,004, Oct. 15, 2015 (6 pages).

United States Patent and Trademark Office, "Non-Final Office Action," issued in connection with U.S. Appl. No. 14/885,216, mailed Mar. 1, 2017 (12 pages).

United States Patent and Trademark Office, "Notice of Allowance," issued in connection with U.S. Appl. No. 14/862,508, on Dec. 28, 2016 (11 pages).

United States Patent and Trademark Office, "Notice of Allowance," issued in connection with U.S. Appl. No. 14/862,508, on Jan. 27, 2017 (4 pages).

Canadian Intellectual Property Office, "Office Action," dated in connection with Canadian Patent Application No. 2,896,096, on Apr. 20, 2017, 3 pages.

The United States Patent and Trademark Office, "Notice of Allowance," dated in connection with U.S. Appl. No. 14/885,216, filed Jul. 7, 2017, 19 pages.

\* cited by examiner

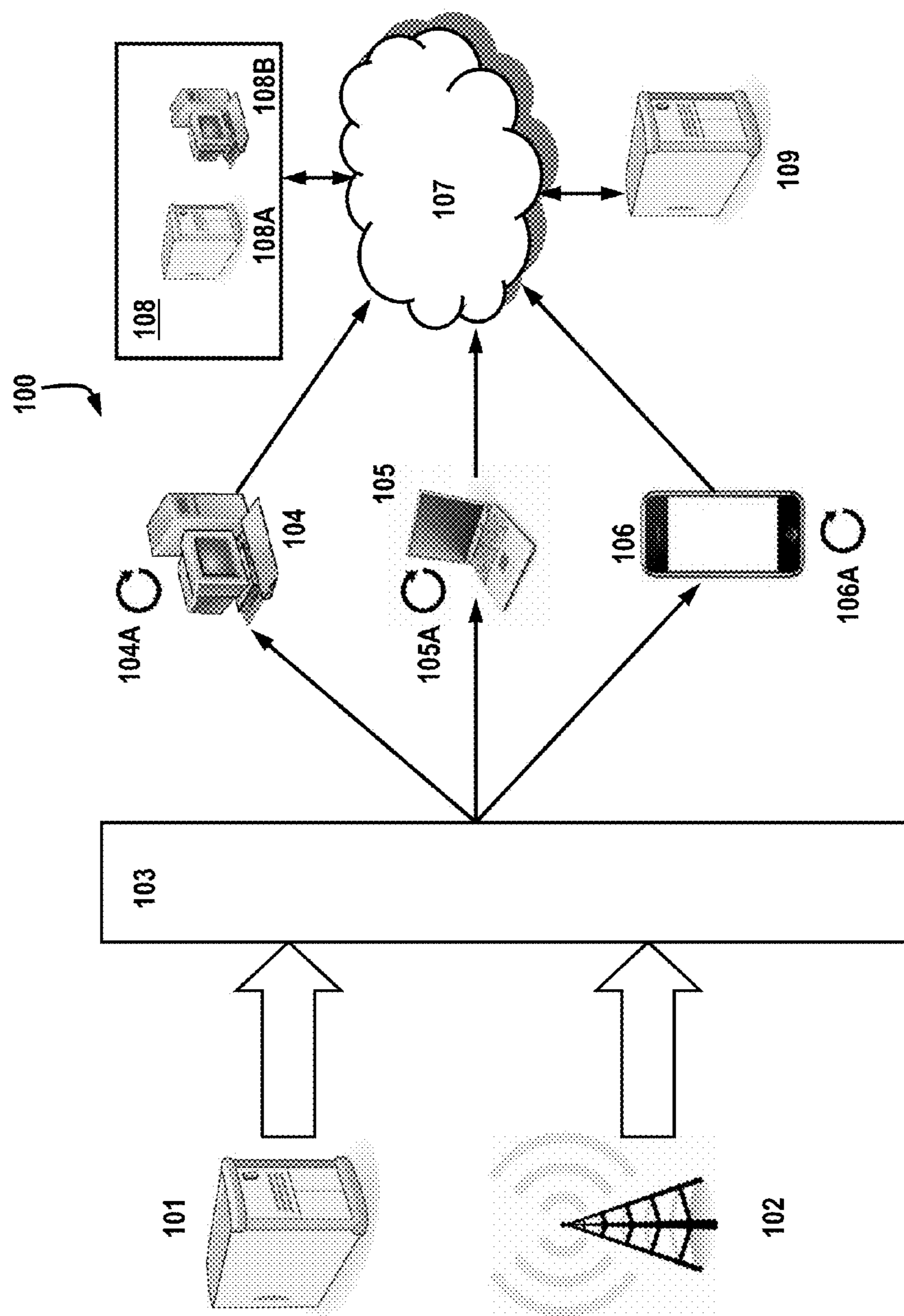


FIG. 1

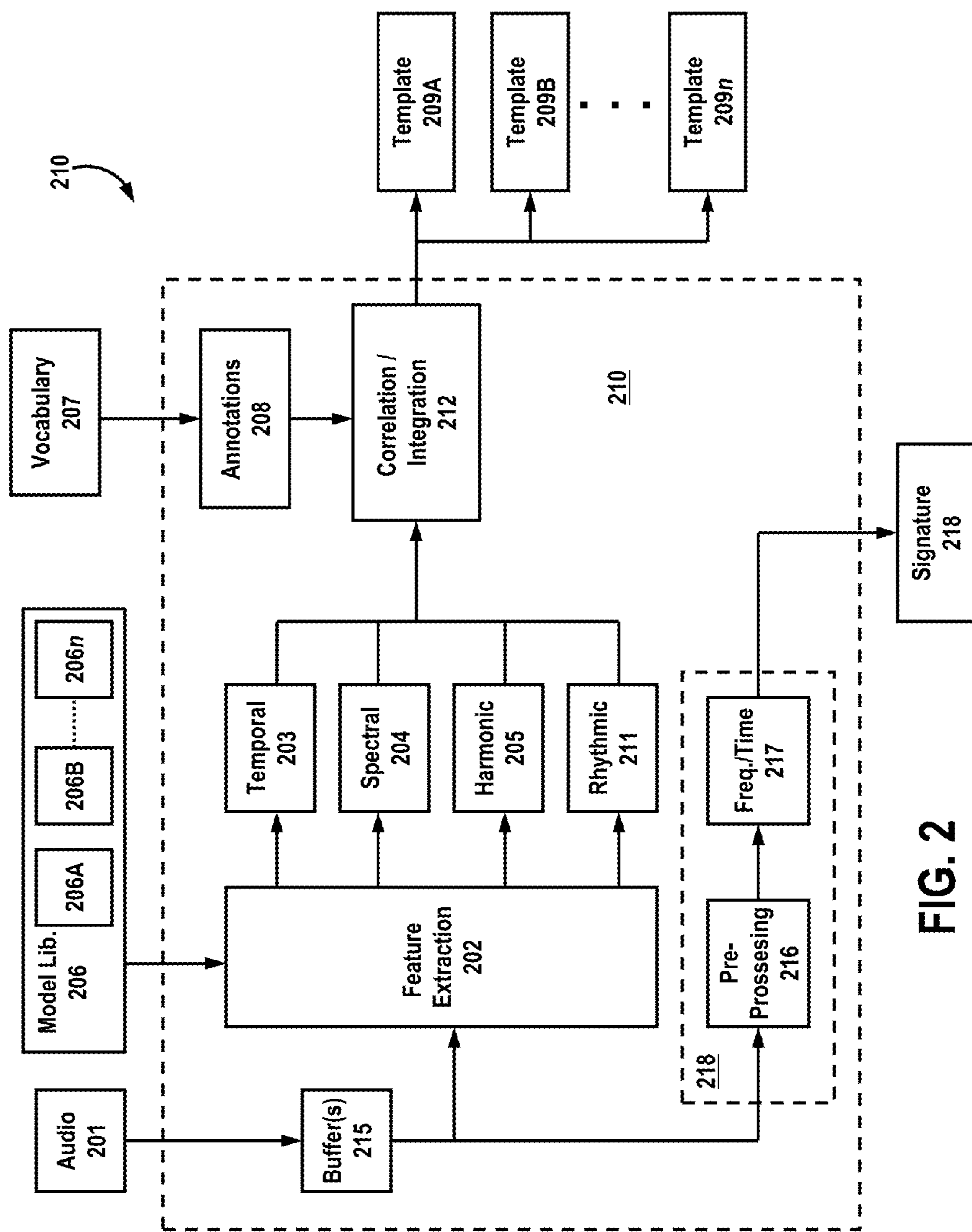


FIG. 2



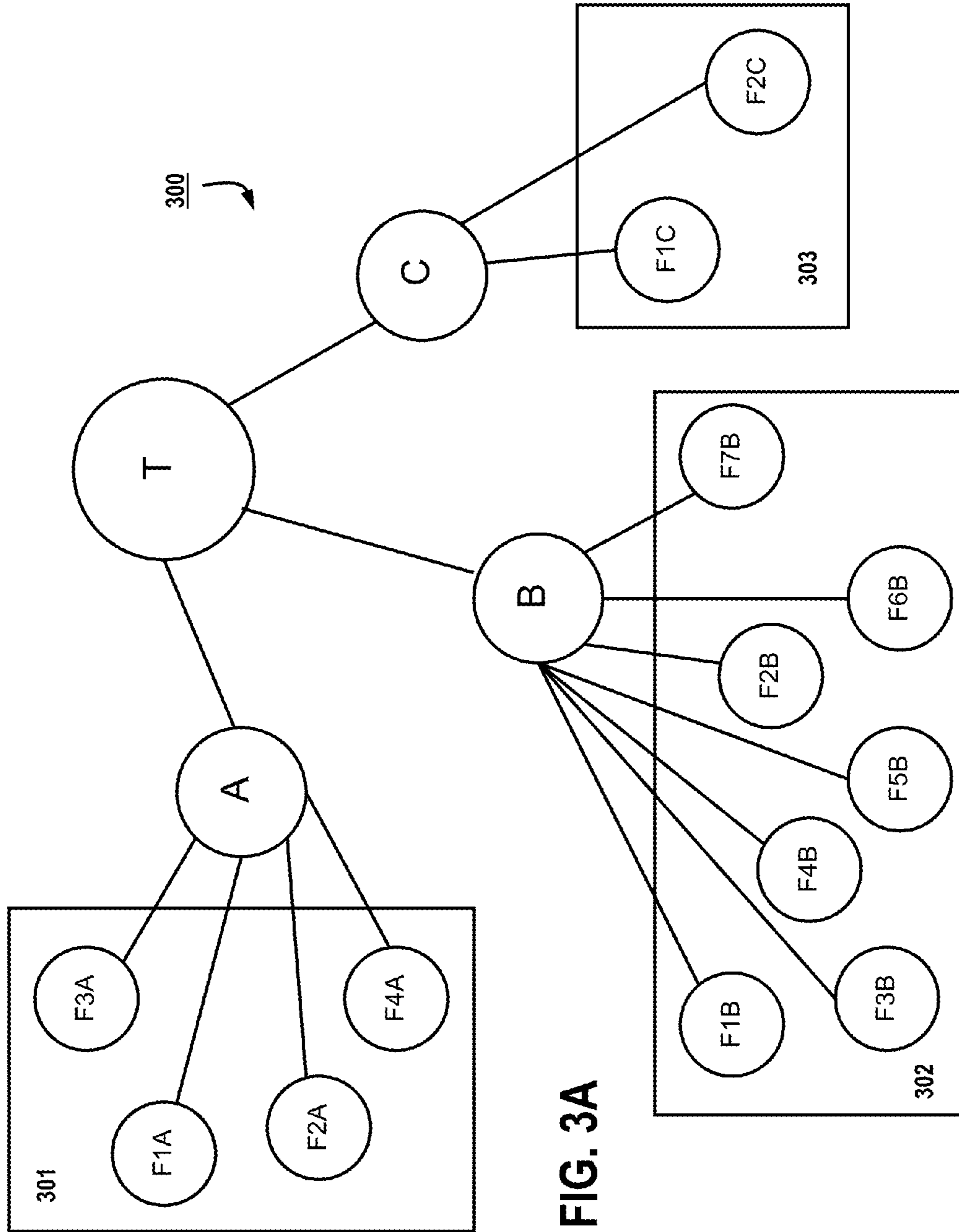


FIG. 3A



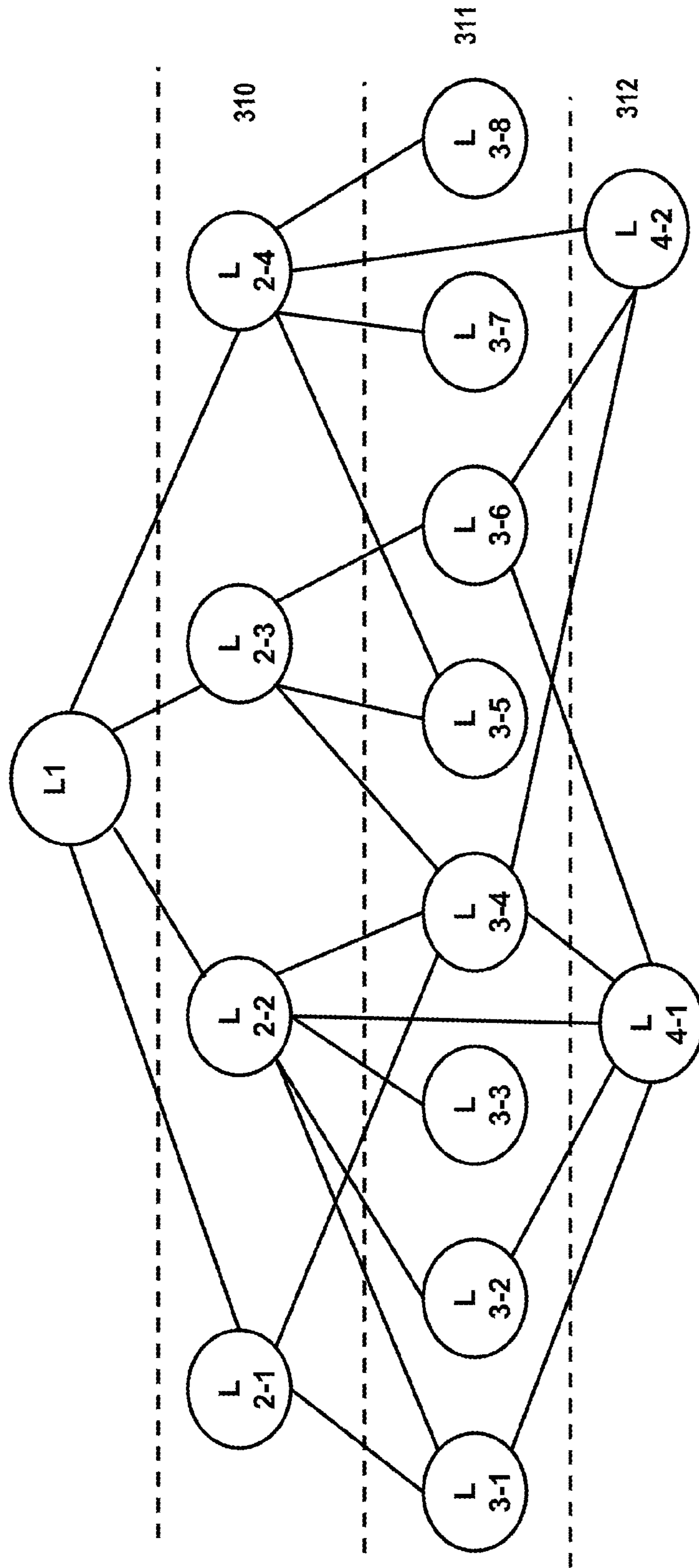


FIG. 3B



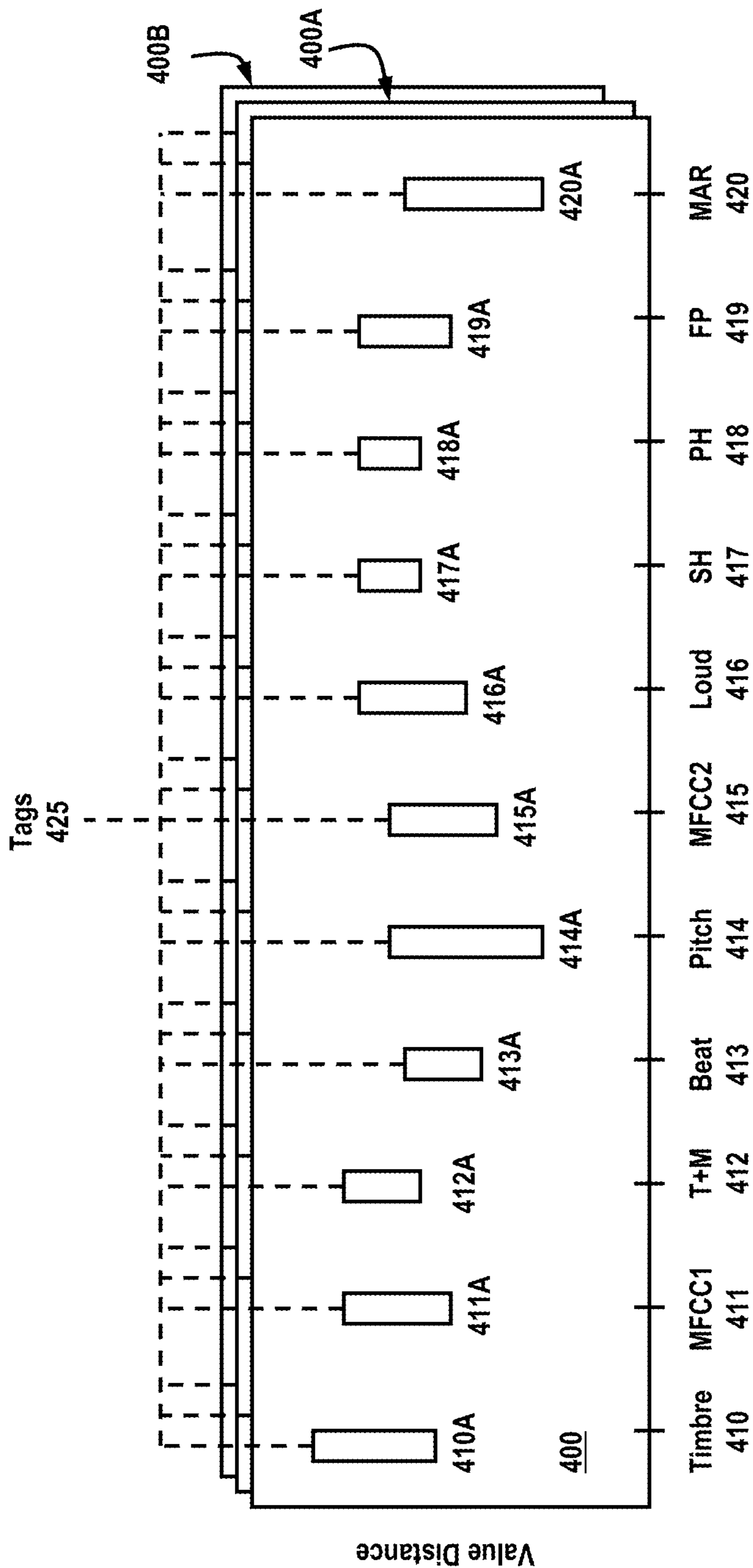
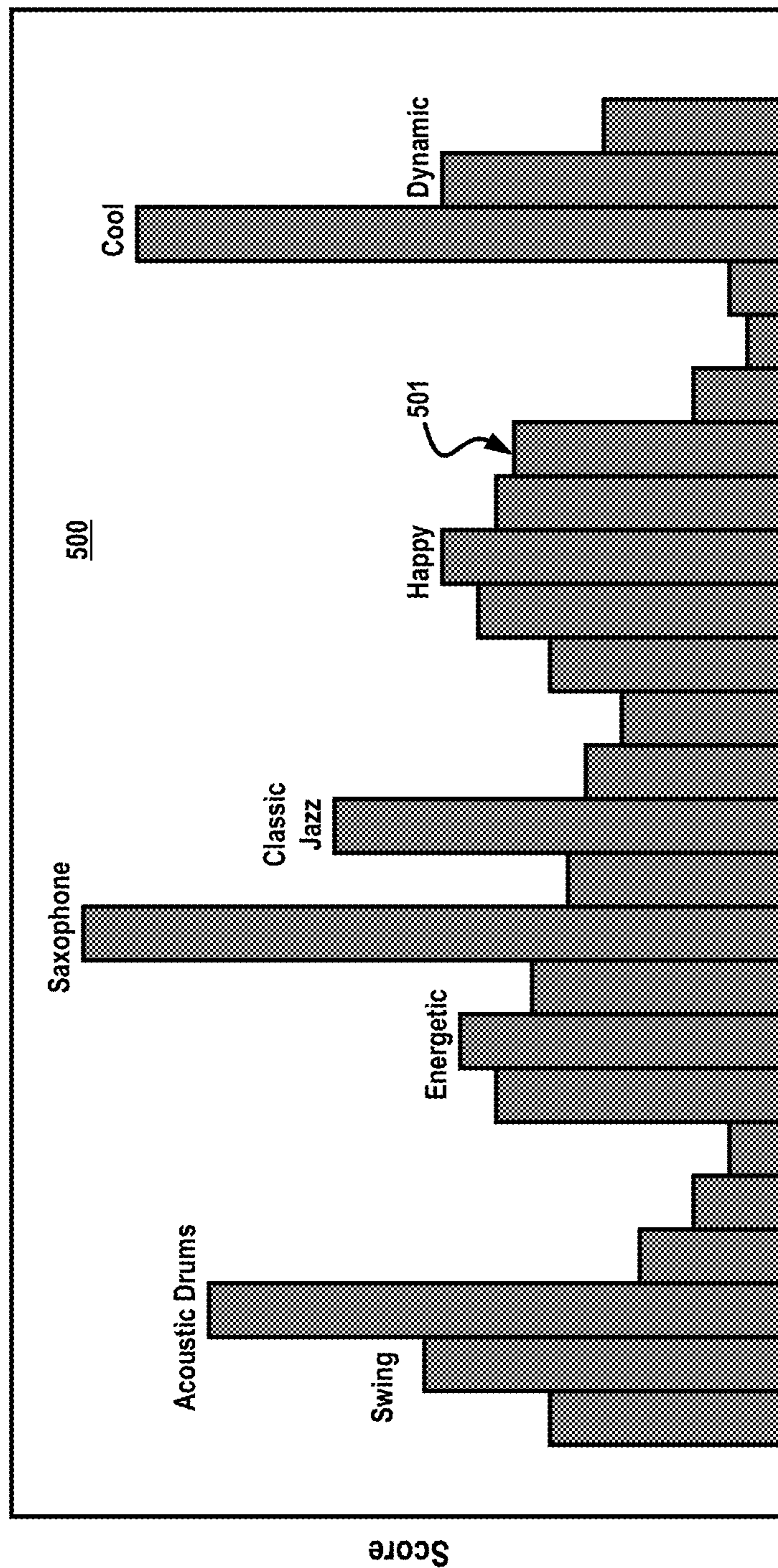


FIG. 4



Tags

FIG. 5

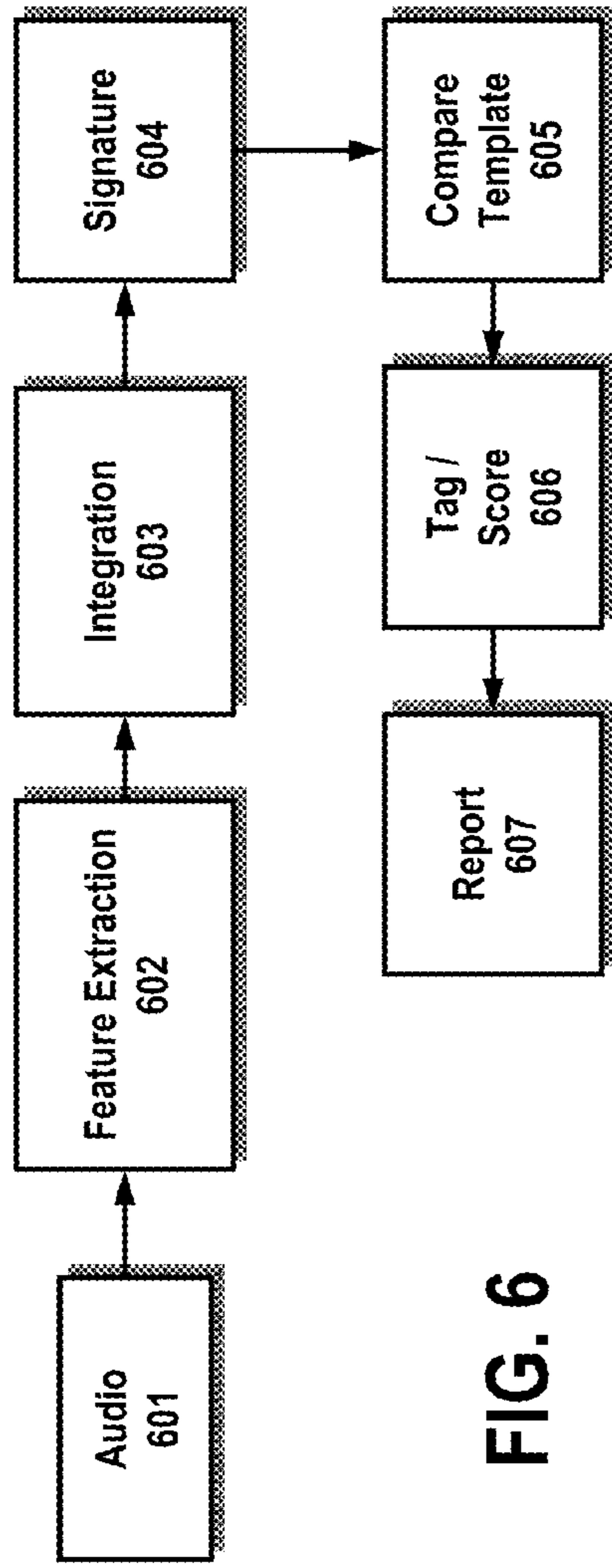


FIG. 6

User00001	
09:20AM – Talk, soothing, conversational	<u>701</u>
09:26AM -- Classical, Baroque, Woodwinds, Happy	
10:20AM -- Jazz, Cool, Saxophone, Swing	
11:20AM -- Classic Rock, Electric Guitar, Brooding	
12:00PM – Talk, confrontational, aggressive	
02:00PM – Hard Rock, energetic, distortion, aggressive	

FIG. 7B

User00001	
09:20AM -- Talk	<u>700</u>
09:26AM -- Classical	
10:20AM -- Jazz	
11:20AM -- Classic Rock	
12:00PM -- Talk	
02:00PM – Hard Rock	

FIG. 7A



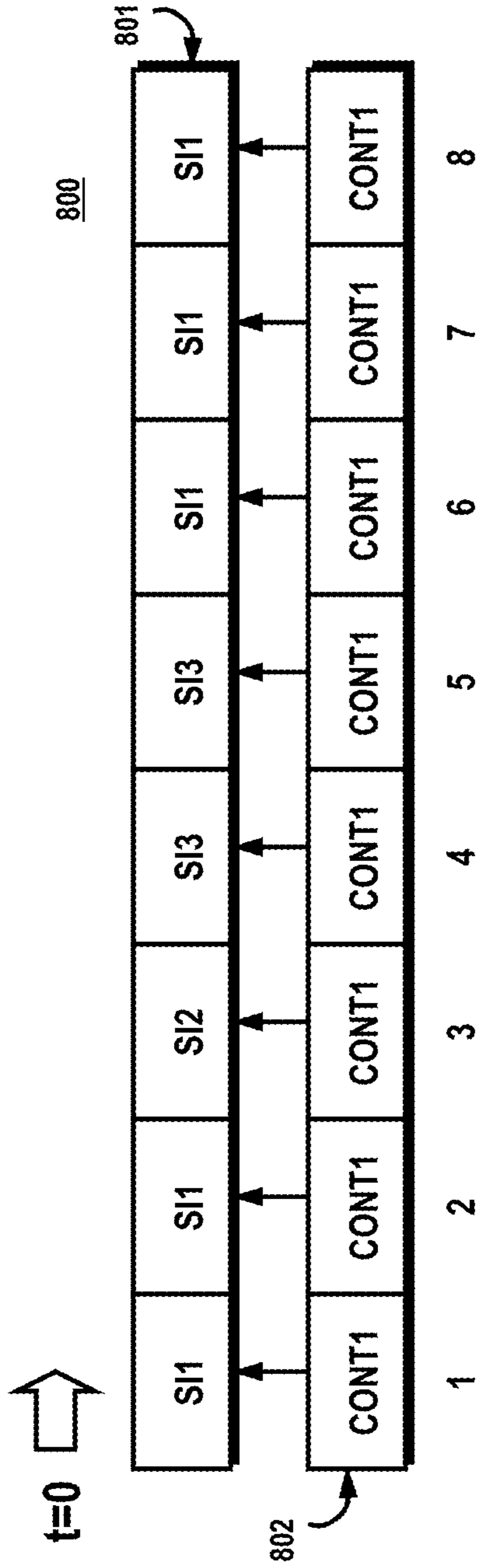


FIG. 8

Time	Signature ID	Semantic Info
901	902	903
1	SIG0035	Info1
2	X	Info1
3	SIG0276	Info6

Diagram 900 is a table with three columns: Time, Signature ID, and Semantic Info. The rows are: 901 (Time), 902 (Signature ID), 903 (Semantic Info); 1 (Time), SIG0035 (Signature ID), Info1 (Semantic Info); 2 (Time), X (Signature ID), Info1 (Semantic Info); 3 (Time), SIG0276 (Signature ID), Info6 (Semantic Info). Arrows connect SIG0035 to X and X to Info1.

FIG. 9

# AUDIO MATCHING WITH SEMANTIC AUDIO RECOGNITION AND REPORT GENERATION

## RELATED APPLICATIONS

The present application is a continuation of U.S. patent application Ser. No. 13/725,004 titled "AUDIO MATCHING WITH SEMANTIC AUDIO RECOGNITION AND REPORT GENERATION" to Neuhauser et al., filed Dec. 21, 2012, which is a continuation-in-part U.S. patent application Ser. No. 13/724,836 titled "Audio Processing Techniques for Semantic Audio Recognition and Report Generation" to Neuhauser et al., filed Dec. 21, 2012, the contents of which are incorporated by reference in entirety herein.

## TECHNICAL FIELD

The present disclosure relates to systems, apparatuses and processes for processing and communicating data, and, more specifically, to process audio portions of media data to extract audio signatures together with semantic audio features, and converting the signatures and features for audience measurement research.

## BACKGROUND INFORMATION

The use of audio "fingerprints" or "signatures" has been known in the art, and was partly pioneered by such companies as Arbitron for audience measurement research. Audio signatures are typically formed by sampling and converting audio from a time domain to a frequency domain, and then using predetermined features from the frequency domain to form the signature. The frequency-domain audio may then be used to extract a signature therefrom, i.e., data expressing information inherent to an audio signal, for use in identifying the audio signal or obtaining other information concerning the audio signal (such as a source or distribution path thereof). Suitable techniques for extracting signatures include those disclosed in U.S. Pat. No. 5,612,729 to Ellis, et al. and in U.S. Pat. No. 4,739,398 to Thomas, et al., both of which are incorporated herein by reference in their entireties. Still other suitable techniques are the subject of U.S. Pat. No. 2,662,168 to Scherbatskoy, U.S. Pat. No. 3,919,479 to Moon, et al., U.S. Pat. No. 4,697,209 to Kiewit, et al., U.S. Pat. No. 4,677,466 to Lert, et al., U.S. Pat. No. 5,512,933 to Wheatley, et al., U.S. Pat. No. 4,955,070 to Welsh, et al., U.S. Pat. No. 4,918,730 to Schulze, U.S. Pat. No. 4,843,562 to Kenyon, et al., U.S. Pat. No. 4,450,551 to Kenyon, et al., U.S. Pat. No. 4,230,990 to Lert, et al., U.S. Pat. No. 5,594,934 to Lu, et al., European Published Patent Application EP 0887958 to Bichsel, PCT Publication WO/2002/11123 to Wang, et al. and PCT publication WO/2003/091990 to Wang, et al., all of which are incorporated herein by reference in their entireties. The signature extraction may serve to identify and determine media exposure for the user of a device.

While audio signatures have proven to be effective at determining exposures to specific media, audio signature systems provide little to no semantic information regarding the media. As used herein below, the terms "semantic," "semantic information," "semantic audio signatures," and "semantic characteristics" refer to information processed from time, frequency and/or amplitude components of media audio, where these components may serve to provide generalized information regarding characteristics of the media, such as genre, instruments used, style, etc., as well as

emotionally-related information that may be defined by a customizable vocabulary relating to audio component features (e.g., happy, melancholy, aggressive). This may be distinguished from "audio signatures" that are used to provide specific information that is used for media content identification, media content distributor identification and media content broadcaster identification (e.g., name of program, song, artist, performer, broadcaster, content provider, etc.).

Some efforts have been made to semantically classify, characterize, and match music genres and are described in U.S. Pat. No. 7,003,515, titled "Consumer Item Matching Method and System," issued Feb. 21, 2006 and is incorporated by reference herein. However, these efforts often rely on humans to physically characterize music. Importantly, such techniques do not fully take advantage of audio signature information together with semantic information when analyzing audio content. Other efforts have been made to automatically label audio content for Music Information Retrieval Systems (MIR), such as those described in U.S. patent application Ser. No. 12/892,843, titled "Automatic labeling and Control of Audio Algorithms by Audio Recognition," filed Sep. 28, 2010, which is incorporated by reference in its entirety herein. However such systems can be unduly complex and also do not take full advantage of audio signature technology and semantic processing. As such, there is a need in the art to provide semantic information based on generic templates that may be used to identify semantic characteristics of audio, and to use the semantic characteristics in conjunction with audio signature technology. Additionally, there is a need to identify such characteristics for the purposes of audience measurement. Currently advertisers target listeners by using radio ratings. These rating are gathered by using encoding or audio matching systems. As listening/radio goes to a one-to-one experience (e.g. Pandora, Spotify, Songza, etc.), there is a need for advertisers to be able to target listeners by the style of music they listen, along with other related information. Semantic analysis can identify this information and provide useful tools for targeted advertisement. Furthermore, semantic information may be used to provide supplemental data to matched audio signature data.

## SUMMARY

Accordingly, under one embodiment, a processor-based method is disclosed for producing supplemental information for audio signature data. The method comprises the steps of receiving the audio signature data at an input from a data network, the audio signature data being received from a device during a first time period, wherein the audio signature data comprises data relating to at least one of time and frequency representing a first characteristic of media content containing audio. Semantic audio signature data may be received at the input (or a different input) from the data network, the semantic audio signature data being received from the device for the first time period, wherein the semantic audio signature comprises at least one of temporal, spectral, harmonic and rhythmic features relating to a second characteristic of the media content, The method then successively associates the semantic audio signature data to the audio signature data in a processor for the first time period.

In another embodiment, a system is disclosed for producing supplemental information for audio signature data. The system comprises an input configured to receive the audio signature data from a data network, the audio signature data being received from a device during a first time period,



wherein the audio signature data comprises data relating to at least one of time and frequency representing a first characteristic of media content containing audio. The input may further be configured to receive semantic audio signature data from the data network, the semantic audio signature data being received from the device for the first time period, wherein the semantic audio signature comprises at least one of temporal, spectral, harmonic and rhythmic features relating to a second characteristic of the media content. A processor may be operatively coupled to the input, where the processor is configured to successively associate the semantic audio signature data to the audio signature data in a processor for the first time period

In yet another embodiment, a processor-based method is disclosed for producing supplemental information for audio signature data. The method comprises the steps of receiving the audio signature data at an input from a data network, the audio signature data being received from a device during a first time period, wherein the audio signature data comprises data relating to at least one of time and frequency representing a first characteristic of media content containing audio. Further, semantic audio signature data may be received at the input from the data network, the semantic audio signature data being received from the device for the first time period, wherein the semantic audio signature comprises at least one of temporal, spectral, harmonic and rhythmic features relating to a second characteristic of the media content. The semantic audio signature data may then be successively associated to the audio signature data in a processor for the first time period, where the associated semantic audio signature data and audio signature data is processed to determine changing second characteristics in relation to a first characteristic.

Additional features and advantages of the various aspects of the present disclosure will become apparent from the following description of the preferred embodiments, which description should be taken in conjunction with the accompanying drawings.

#### BRIEF DESCRIPTION OF THE DRAWINGS

FIG. 1 is a block diagram illustrating a media measurement system under one exemplary embodiment;

FIG. 2 illustrates one configuration for generating audio templates and signatures for use in extracting semantic features from audio under an exemplary embodiment;

FIG. 2A illustrates an exemplary signature formation process under one embodiment;

FIG. 3A illustrates audio feature template arrangement under another exemplary embodiment;

FIG. 3B illustrates an audio feature template hierarchy under another exemplary embodiment;

FIG. 4 illustrates an exemplary process for generating tags for use in audio template generation under yet another exemplary embodiment;

FIG. 5 illustrates an exemplary process for processing audio samples for comparison with audio templates to provide tag scores under yet another exemplary embodiment;

FIG. 6 illustrates an exemplary tag score utilizing the audio processing described above in FIG. 5;

FIGS. 7A and 7B illustrate exemplary reports that may be generated from tag scoring under another exemplary embodiment;

FIG. 8 illustrates an exemplary embodiment where audio signatures are combined with semantic information to represent the semantic development of content; and

FIG. 9 illustrates an exemplary embodiment, where semantic information is used to supplement audio signature information.

#### DETAILED DESCRIPTION

FIG. 1 is an exemplary block diagram for a system 100, wherein media is provided from a broadcast source 102 (e.g., television, radio, etc.) and/or a data source 101 (e.g., server, cloud, etc.). The media is communicated to a media distribution network 103, which has the ability to pass through the broadcast and/or data to remote users or subscribers. Such media distribution networks 103 are well known and may include broadcast stations, satellite/cable, routers, servers, and the like.

The media may be received at one or more locations using any of a number of devices, including a personal computer 104, laptop 105, and smart phone or tablet 106. It is understood by those skilled in the art that the present disclosure is not limited strictly to devices 104-106, but may include any device configured to receive and/or record media including set-top-boxes, IPTV boxes, personal people meters, and the like. Additionally, devices, such as 104-106 may be equipped with one or more microphones (not shown) for transducing ambient audio for sampling and processing. Examples of such configurations may be found in U.S. patent application Ser. No. 13/341,272, titled "Apparatus, System And Method For Activating Functions In Processing Devices Using Encoded Audio And Audio Signatures," filed Dec. 30, 2011, and U.S. patent Ser. No. 13/341,365, titled "Activating Functions In Processing Devices Using Start Codes Embedded In Audio" filed Dec. 30, 2011, both of which are incorporated by reference in their entireties herein. Devices 104-106 may also be capable of reproducing media (104A-106A) on the device itself, where the media is transferred, downloaded, stored and/or streamed.

As each device 104-106 receives media from network 103 and/or reproduces media locally 104A-106A, the audio portion of the media is sampled and processed to form semantic audio signatures or templates, where resulting signature data is time stamped and transmitted to computer network 107 via wired or wireless means that are known in the art. In addition to semantic audio signature data, devices 104-106 may additionally transmit identification information that identifies the device and/or the user registered for the device. Under one embodiment, demographic information relating to the users of any of devices 104-106 may be transmitted as well. The semantic signatures are then stored in one or more remote locations or servers 109, where they are compared with audio signature templates provided from system 108 for semantic audio analysis. Under one exemplary embodiment, system 108 comprises at least one workstation 108B and server 108A, where audio signature templates are produced using any of the techniques described below, and forwarded to server(s) 109.

Turning to FIG. 2, an exemplary configuration is shown for creating audio signature templates, where audio is fed into a processing device 210. It is understood that processing device 210 of FIG. 2 may be a dedicated workstation (e.g., 108B), or a portable devices, such as a smart phone, tablet, PC, etc. (104A-106A). Under one exemplary embodiment, audio 201 is sampled and stored in one or more buffers (215), where portions of the audio are processed and subjected to one or more feature extractions (202). Additionally, music portions stored in buffers 215 are subjected to signature extraction, which will be discussed in greater detail below.



## 5

With regard to feature extraction **202**, while specific feature sets may vary, under a preferred embodiment, extracted feature sets in **204** may include energy-based features, spectral features, rhythmic features, temporal features and/or harmonic features. Depending on the feature set used, different models (**206A-206n**) may be called from a model library **206** memory in order to facilitate appropriate feature extraction. The feature extraction process is preferably controlled by software operative on a tangible medium, such as Psysound (<http://psysound.wikidot.com/>), CLAM (<http://clam-project.org/>), Marsyas (<http://marsyas.sness.net/>), MIRToolbox (<https://www.jyu.fi/hum/laitokset/musiikki/en/research/coe/materials/mirtoolbox>), MA Toolbox (<http://www.ofai.at/~elias.pampalk/ma/>), Sound Description Toolbox, and/or any other suitable program or application, preferably compatible with the MATLAB and MPEG-7 format.

Feature extraction in **202** may advantageously be separated into multiple stages, where, for example, a first stage is responsible for processing temporal features **203**, while a second stage is independently responsible for processing spectral features **204**. Alternately or in addition, the stages may be separated by sample size, so that longer samples are processed for certain features before shorter sub-samples are processed for other features. This configuration may be advantageous for extracting features that are optimally detected over longer periods of time (e.g., 30 sec.), while reserving shorter segments (e.g., 5-6 sec., 100-200 ms) for other feature extraction processes. The varying sample sizes are also useful for separating audio segments that are independently processed for audio signature extraction **218**, since audio signature extraction may rely on audio portions that are smaller than those required for certain templates.

For semantic processing, feature extraction **202** preferably includes pre-processing steps such as filtering and normalization to provide zero mean and unity variance. A first-order finite impulse response (FIR) filter may also be used to increase the relative energy of high-frequency spectrum. Frame blocking or "windowing" is then performed to segment the signal into statistically stationary blocks. The frame size (in terms of sample points) should be equal to the powers of 2 (such as 256, 512, 1024, etc) in order to make it suitable for transformation (e.g., FFT). Hamming window may be used to weight the pre-processed frames. In order to reduce the difference between neighboring frames, an overlap may be applied that is up to  $\frac{2}{3}$  of the original frame size. However, the greater the overlap, the more computational power is needed.

For temporal feature extraction **203**, features are taken from the native domain of the audio signal, and without any preceding transformation. Temporal features include, but are not limited to, amplitude, power, and zero-crossing of the audio signal. Amplitude based features are processed directly from the amplitudes of a signal and represent the temporal envelope of an audio signal. Utilizing an audio waveform descriptor (e.g., MPEG-7), a compact description of the shape of a waveform may be formed by computing the minimum and maximum samples within non-overlapping portions of frames, resulting in a representation of the (preferably down-sampled) waveform envelope over time. Also, amplitude descriptors may be used by separating the audio signal into segments having low and high amplitudes according to an adaptive threshold. The duration, variation of duration and energy of segments crossing the thresholds would be recorded to form a specific descriptor for an audio segment. The amplitude descriptor could thus be used to

## 6

characterize audio in terms of quiet and loud segments and to distinguish audio with characteristic waveform envelopes.

Generally speaking, the energy of a signal is the square of the amplitude of a waveform, and power may be represented are the transmitted energy of the signal per unit of time. Short Time Energy (STE) processing may be performed on the envelope of a signal to determine mean energy per frame. Thus, power may be represented as the mean square of a signal. Root-Mean-Square (RMS) may used to measure the power (or loudness, volume) of a signal over a frame. The global energy of a signal  $x$  can be computed by taking the root average of the square of the amplitude (RMS), expressed by

$$x_{RMS} = \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2}$$

Additionally, a temporal centroid (MPEG-7) may be used to determine a time average over the envelope of a signal to determine a point(s) in time where most of the energy of the signal is located on average. Such features are advantageous for distinguishing percussive from sustained sounds.

Continuing with the temporal features, the zero crossing rate (ZCR) may be used to measure signal noisiness and may be calculated by taking the mean and standard deviation of the number of signal values that cross the zero axis in each time window (i.e., sign changes of the waveform:

$$ZCR = \frac{1}{T} \sum_{t=m-T+1}^m \frac{|\text{sgn}(s_t) - \text{sgn}(s_{t-1})|}{2} w(m-t)$$

where  $T$  is the length of a time window,  $s_t$  is the magnitude of the  $t$ -th time-domain sample and  $w$  is a rectangular window. The ZCR is advantageous in discriminating between noise, speech and music, where the ZCR would be greatest for noise, less for music, and lesser still for speech. Additional techniques, such as linear prediction zero crossing ratios could be used to determine a ratio of the zero crossing count of a waveform and the zero crossing count of the output of a linear prediction analysis filter. Such a feature would be advantageous in determining the degree of correlation in a signal.

While time domain features may provide useful data sets for semantic audio analysis, even more valuable information may be obtained from the spectral domain. To achieve this, a transformation should be performed on an audio signal to convert time domain features to the spectral domain wherein the existence and progression of periodic elements may be obtained, as well as pitch, frequency ranges, harmonics, etc. The most common and well-known transformation is the Fourier Transformation. For discrete digital signals, a Discrete Fourier Transformation (DFT) is generally known as

$$X_k = \sum_{n=0}^{N-1} x_n e^{-\frac{2\pi i}{N} nk} \quad k = 0, 1, \dots, N-1$$

where  $x_n$  are the time-based complex numbers,  $N$  is the number of values to transform and  $X_k$  the resulting Fourier transformed complex numbers (or "Fourier coefficients").

The spectral domain (**204**) allows several extractions and computational time-invariant possibilities that bring out



characteristic features and representations like spectrograms, energy deviations, frequency histograms and magnitudes of certain frequency range transformation that illustrate its influence on human perception of audio. In addition, time discrete Short Time Fourier Transformations (STFT) are preferably performed on short single segments of audio that is changing over time, resulting in a representation of frequency lots at a specific time, which may further be depicted in a time-frequency plot that may further be semantically processed using Bark scales. The Bark scale is a psycho acoustic scale that matches frequency range intervals to a specific number, and is based on the perception of pitch for human beings with respect to the amount of acoustic “feeling.” It considers the almost-linear relation in lower frequency ranges as well as the logarithmic in higher ranges and its basic idea originates from frequency grouping and the “subdivision concept” referred to in the area of human hearing. As STFT may produce real and complex values, the real values may be used to process the distribution of the frequency components (i.e., spectral envelope) while the complex values may be used to process data relating to the phase of those components.

Under one embodiment, spectral features **204** are extracted under STFT, and, depending on the model used, may produce timbral texture features including spectral centroid, spectral rolloff, spectral flux, spectral flatness measures (SFM) and spectral crest factors (SCF). Such features are preferably extracted for each frame and then by taking the mean and standard deviation for each second. The sequence of feature vectors may be combined and/or collapsed into one or more vectors representing the entire signal by taking again the mean and standard deviation. A spectral centroid (SC) refers to the centroid, or “center of gravity” of the magnitude spectrum of the STFT and may be expressed as

$$SC = \frac{\sum_{n=1}^N nA_t^n}{\sum_{n=1}^N A_t^n}$$

where  $A_t^n$  is the magnitude of the spectrum at the t-th frame and the n-th frequency bin, and N is the total number of bins. As the centroid provides a measure of spectral shape, a higher spectral centroid will indicate a “brighter” audio texture.

The spectral rolloff is a spectral feature that estimates the amount of high frequency in a signal. More specifically, spectral rolloff may be defined as the frequency  $k_t$  below which a certain fraction or percentage of total energy is contained. This fraction may be fixed by default to a specific number, such as 0.85 or 0.95, such as:

$$\sum_{n=1}^{k_t} A_t^n = 0.85 * \sum_{n=1}^{k_t} A_t^n$$

Spectral Flux (SF) estimates the amount of local spectral change and may be defined as a spectral feature representing the square of the difference between the normalized magnitudes of successive frames:

$$SF = \sum_{n=1}^N (a_t^n - a_{t-1}^n)^2$$

where  $a$  denotes the normalized magnitude of the spectrum that is preferably normalized for each frame. Because spectral flux represents the spectral variations between adjacent frames, it may be correlated to features such as articulation.

Tonalness is an audio feature that is useful in quantifying emotional perceptions, where joyful or peaceful melodies may be characterized as being more tonal (tone-like), while angry or aggressive melodies are characterized as being more atonal (noise-like). Features indicating tonalness include spectral flatness (SFL) and spectral crest factors (SCF), where SFL is expressed as the ratio between the geometric mean of the power spectrum and its arithmetic mean:

$$SFL = \frac{\left(\prod_{n \in B^k} A_t^n\right)^{1/N_k}}{\frac{1}{N_k} \sum_{n \in B^k} A_t^n}$$

and SCF is the ratio between the peak amplitude and the RMS amplitude:

$$SCF = \frac{\max_{n \in B^k} A_t^n}{\frac{1}{N_k} \sum_{n=1}^N A_t^n}$$

where  $B^k$  denotes the k-th frequency subband and  $N_k$  is the number of bins in  $B^k$ . While any suitable number of subbands may be used, under one exemplary embodiment, 24 subbands are used for SFL and SCF extraction.

In addition to the above spectral features, it may be advantageous to also extract Mel-frequency cepstral coefficients (MFCCs) which denote coefficients of a discrete cosine transform (DCT) of each short-term log power spectrum expressed on a non-linear perceptually-related Mel-frequency scale to represent the formant peaks of the spectrum. Preferably, the mean and standard deviation of the first 13-20 MFCCs of each frame are taken. Since MFCCs average the spectral distribution in each subband, some of the relative spectral information may be lost. To compensate for this, and octave based spectral contrast (OBSC) may be used to capture the relative energy distribution of the harmonic components in the spectrum. OBSC considers the strength of spectral peaks and valleys in each sub-band separately. In general, spectral peaks correspond to harmonic components and spectral valleys correspond to non-harmonic components or noise in a music piece. Therefore, the contrast between spectral peaks and spectral valleys will reflect the spectral distribution.

Under another embodiment, spectral features may include the extraction of Daubechies wavelets coefficient histograms (DWCH), which is computed from the histograms of Daubechies wavelet coefficients at different frequency subbands with different resolutions, and is described in U.S. patent application Ser. No. 10/777,222, titled “Music Feature Extraction Using Wavelet Coefficient Histograms”, filed Feb. 13, 2004, and is incorporated by reference in its entirety herein.



As different spectral features are extracted, they may be used or combined to form sensory features such as spectral dissonance, irregularity and inharmonicity. Spectral dissonance measures the noisiness of the spectrum, where notes that do not fall within a prevailing harmony are considered dissonant. Spectral dissonance may be estimated by computing the peaks of the spectrum and taking the average of all the dissonance between all possible pairs of peaks.

Irregularity measures the degree of variation of the successive peaks of the spectrum and may be computed by summing the square of the difference in amplitude between adjoining partials, or

$$\frac{\sum_{n=1}^N (A_t^n - A_t^{n+1})^2}{\sum_{n=1}^N A_t^n * A_t^{n+1}}$$

Alternately, irregularity may be measured using Krimphoff's method, which defines irregularity as the sum of amplitude minus the mean of the preceding, current, and next, amplitude:

$$\sum_{n=2}^{N-1} \left| A_t^n - \frac{A_t^{n-1} + A_t^n + A_t^{n+1}}{3} \right|$$

Inharmonicity estimates the amount of partials that depart from multiples of the fundamental frequency. It is computed as an energy weighted divergence of the spectral components from the multiple of the fundamental frequency, or

$$\frac{2}{f_0} \frac{\sum_{n=1}^N |f_n - nf_0| (A_t^n)^2}{\sum_{n=1}^N (A_t^n)^2}$$

where  $f_n$  is the n-th harmonic of the fundamental frequency  $f_0$ . The inharmonicity represents the divergence of the signal spectral components from a purely harmonic signal. The resulting value ranges from 0 (purely harmonic) to 1 (inharmonic)

Still referring to FIG. 2, harmonic feature extraction may also be performed to extract features from the sinusoidal harmonic modeling of an audio signal. Harmonic modeling may be particularly advantageous for semantic analysis as natural/musical sounds are themselves harmonic, consisting of a series of frequencies at multiple ratios of the lowest frequency, or fundamental frequency  $f_0$ . Under one embodiment, a plurality of pitch features (e.g., salient pitch, chromagram center) and tonality features (e.g., key clarity, mode, harmonic change) are extracted. The perceived fundamental frequency of a time frame (e.g., 50 ms, 50% overlap) may be calculated using a multi-pitch detection algorithm by decomposing an audio waveform into a plurality of frequency bands (e.g., one below and one above 1 kHz), computing an autocorrelation function of the envelope in each subband, and producing pitch estimates by selecting the peaks from the sum of the plurality of autocorrelation functions. The calculation corresponding to the highest peak is deemed the "salient pitch."

Alternately or in addition, a pitch class profile or wrapped chromagram may be computed for each frame (e.g., 100 ms, 1/8 overlap), where the centroid of the chromagram is selected as the fundamental frequency, or chromagram centroid. A wrapped chromagram may project a frequency spectrum onto 12 bins representing 12 semitones (or chroma) of a musical octave (e.g., 440 Hz (C4) and 880 Hz (C5) would be mapped to chroma "C"). By comparing a chromagram to the 24 major and minor key profiles, key detection may be performed to estimate the strength of a frame compared to each key (e.g., C major). The key associated with the greatest strength would be identified as the key strength or key clarity. The difference between the best major key and best minor key in strength may be used as an estimate of music mode, which may be used to characterize a fixed arrangement of the diatonic tones of an octave. The numerical value would be indicative of audio content being more major, and thus having a higher value.

Harmonic changes (e.g., chord changes) may also be determined using a Harmonic Change Detection Function (HCDF) algorithm modeled for equal tempered pitch space for projecting collections of pitches as tonal centroid points in a 6-D space. The HCDF system comprises a constant-Q spectral analysis at the lowest level, followed by a 12-semitone chromagram decomposition. A harmonic centroid transform is then applied to the chroma vectors which is then smoothed with a Gaussian filter before a distance measure is calculated. High harmonic change would indicate large differences in harmonic content between consecutive frames. Short term features could be aggregated by taking mean and standard deviation. Additional information on HCDF techniques may be found in Harte et al., "Detecting Harmonic Changes in Musical Audio," AMCOMM '06 Proceedings of the 1st ACM workshop on Audio and music computing multimedia, pp. 21-26 (2006).

Alternately or in addition, a pitch histogram may be calculated using Marsyas toolbox, where pluralities of features may be extracted from it, including tonic, main pitch class, octave range of dominant pitch, main tonal interval relation, and overall pitch strength. Modules such as Pysound may be used to compare multiple pitch-related features including the mean, standard deviation, skewness and kurtosis of the pitch and pitch strength time series.

Additionally, rhythmic features may be extracted from the audio signal. One beat detector structures may comprise a filter bank decomposition, followed by an envelope extraction step, followed by a periodicity detection algorithm to detect the lag at which the signal's envelope is most similar to itself. The process of automatic beat detection may be thought of as resembling pitch detection with larger periods (approximately 0.5 s to 1.5 s for beat compared to 2 ms to 50 ms for pitch). The calculation of rhythmic features may be based on the wavelet transform (WT), where WT provides high time resolution and low-frequency resolution for high frequencies, and low time and high-frequency resolution for low frequencies. The discrete wavelet transform (DWT) is a special case of the WT that provides a compact representation of the signal in time and frequency that can be computed efficiently using a fast, pyramidal algorithm related to multi-rate filterbanks.

The feature set for representing rhythm structure may be based on detecting the most salient periodicities of the signal. The signal may be first decomposed into a number of octave frequency bands using the DWT. Following this decomposition, the time domain amplitude envelope of each band is extracted separately. This is achieved by applying full-wave rectification, low pass filtering, and down-sam-



pling to each octave frequency band. After mean removal, the envelopes of each band are then summed together and the autocorrelation of the resulting sum envelope is computed. The dominant peaks of the autocorrelation function correspond to the various periodicities of the signal's envelope. These peaks are accumulated over the whole sound file into a beat histogram where each bin corresponds to the peak lag, i.e., the beat period in beats-per-minute (BPM). The amplitude of each peak is preferably added to the beat histogram so that, when the signal is very similar to itself (i.e., strong beat) the histogram peaks will be higher. The beat histogram may be processed to generate additional features, such as beat strength, amplitude and period of the first and second peaks of the beat histogram, and the ratio of the strength of the two peaks in terms of BPMs.

Rhythm patterns may also be extracted by calculating a time-invariant representation for the audio signal to provide data on how strong and fast beats are played within the respective frequency bands. The amplitude modulation of the loudness sensation per critical-band for each audio frame sequence (e.g., 6 sec,) is calculated using a FFT. Amplitude modulation coefficients may be weighted based on the psychoacoustic model of the fluctuation strength. The amplitude modulation of the loudness has different effects on human hearing sensations depending on the modulation frequency. The sensation of fluctuation strength tends to be most intense around 4 Hz and gradually decreases up to a modulation frequency of 15 Hz. For each frequency band, multiple values for modulation frequencies between specific ranges (e.g., 0 and 10 Hz) are obtained to indicate fluctuation strength. To distinguish certain rhythm patterns better and to reduce irrelevant information, gradient and Gaussian filters may be applied. To obtain a single representation for each audio signal **201** input into **210**, the median of the corresponding sequences may be calculated to produce an X by Y matrix. A rhythm pattern may be further integrated into a multi-bin (e.g., 60-bin) rhythm histogram by summing amplitude modulation coefficients across critical bands. The mean of the rhythm histogram may be regarded as an estimate of the average tempo.

Rhythm strength may be calculated as the average onset strength of an onset detection curve using algorithmic processes described in Anssi Klapuri, "Sound Onset Detection by Applying Psychoacoustic Knowledge," Proceedings., 1999 IEEE International Conference on Acoustics, Speech, and Signal Processing, vol. 6, pp. 3089-3092 (1999), where the "onset" refers to the start of each musical event (e.g., note). Rhythm regularity and rhythm clarity may be computed by performing autocorrelation on the onset detection curve. If a music segment has an obvious and regular rhythm, the peaks of the corresponding autocorrelation curve will be obvious and strong as well. Onset frequency, or event density, is calculated as the number of onset notes per second, while tempo may be estimated by detecting periodicity from the onset detection curve.

Still referring to FIG. 2, each of the temporal **203**, spectral **204**, harmonic **205**, and rhythmic **211** features are correlated to the audio **201** in **212** to arrange a base set of features. These features may be define in system **210** using vocabulary database **207** that contains a lexicography of various and different words/phrases used to tag the semantic information contained in **212**. Under a preferred embodiment, vocabulary **207** is customizable by an operator of system **210**, where specific words, phrases and descriptions may be entered, depending on the need and audio features involved. For example, in a very simple configuration, the vocabulary may comprise a few genres, styles, and emotive descriptors,

where descriptive words/phrases (tags) are mapped to respectively extracted features. In one embodiment, descriptive tags may be mapped to multiple extracted features. Such a configuration is advantageous in instances where multiple variations of a specific feature (e.g., beat histogram) may be attributable to a single tag (e.g., genre, emotive descriptor).

In another embodiment, entries in the vocabulary are subjected to an annotation process **208** which is advantageous for creating more complex multiclass, multi-tag arrangements and classifications, where tags are arranged in a class, sub-class hierarchy. A class-conditional distribution may then be formed during a training process to attribute tags to extracted features that are positively associated with that tag. The tags may then be arranged in a rank order for later processing and identification using techniques such as Byes' rule, k-nearest neighbor, and fuzzy classification, among others. Once extracted features are tagged and integrated in **212**, a plurality of templates **209A-209n** are formed comprising one or more tagged extracted features, and stored for later comparison.

FIG. 3A provides one example of a template arrangement **300**, where tag T is comprised of 3 extracted audio features tagged as A-C. In this example, tagged feature A is associated with extracted features F1A-F4A **301**, tagged feature B is associated with features F1B-F7B **302**, and tagged feature C is associated with extracted features F1C-F2C **303**. In each of these (**301-303**), extracted features may be values associated with the temporal **203**, spectral **204**, harmonic **205** and/or rhythmic **211** processing performed in FIG. 2. Under a preferred embodiment, certain individual extracted features **301-33** may be duplicated among the tags (A-C), to simplify the datasets used for a tree hierarchy.

FIG. 3B exemplifies one possible hierarchy arrangement where a global tag L1 represents the overall characteristics of extracted features and is labeled according to an assigned vocabulary. In this example, global tag L1 is characterized by four lower-level (**310-312**) tags (L2-1 through L2-4). Each of these lower-level tags may represent different features as a class that may be extracted from different aspects of audio (e.g., temporal, spectral, harmonic, rhythmic), which may be correlated and cross-correlated as shown in FIG. 3B. Below level **310**, is a first sub-level **311** provides additional features, followed by a second sub-level **312** having further additional features that also are correlated and/or cross-correlated. It is understood by those skilled in the art that the tags and level hierarchies may be arranged in a myriad of ways, depending on the needs of the designer. For example, under an alternate embodiment, global tags may represent any of genre, emotional descriptor, instrument, song style, etc. Mid-level features may be associated with lower-level tags representing rhythmic features, pitch and harmony. A sub-level may include tags representing low-level features such as timbre and temporal features. Tags may had additional annotations associated with their class as well, e.g., rhythm (sub: beat histogram, BPM), pitch (sub: salient pitch, chromagram center), timbre (sub: ZCR, SC, SFL, MFCC, DWCH). Moreover, the hierarchical arrangement may be configured to separately take into consideration short-term audio features (e.g., timbre) and long-term audio features (e.g., temporal, pitch, harmony).

The arrangement may be useful in taking advantage of numerous classification techniques for forming and identifying semantic audio information. Under one embodiment, each audio frame is classified separately, and classification results are combined over an analysis segment to get a global classification result. Under another embodiment, the temporal relationship between frames may be taken into



account. One exemplary classifier is a k-Nearest Neighbor Classifier, where the distance between tested tagged feature vectors and the training vectors is measured, and the classification is identified according to the k nearest training vectors. A Gaussian Mixture Model may be used to obtain distributions of feature values for specific musical characteristics, and may be modeled as a weighted sum of Gaussian density functions. This mixture may be used to determine the probability of a test feature vector as belonging to a particular audio characteristic.

In one embodiment, tree-based vector quantization may be used to model discrimination function between classes defined by a set of labeled codebook vectors. A quantization tree is formed to partition the feature space into regions with maximally different tag/class populations. The tree may be used to form a histogram template for an audio characteristic and the classification may be done by matching template histograms of training data to the histograms of the test data. The classification can alternately be done with a feed-forward neural network that is trained with examples from different classes so as to map the high-dimensional space of feature vectors onto the different classes. A Linear Discriminant Analysis (LDA) may be used to find a linear transformation for the feature vectors that best discriminates them (e.g., using Euclidean distance) among classes.

In an alternate embodiment, a binary classification approach may be done using Support Vector Machines (SVMs), where feature vectors are first non-linearly mapped into a new feature space and a hyperplane is then searched in the new feature space to separate the data points of the classes with a maximum margin. The SVM may be extended into multi-class classification with one-versus-the-rest, pairwise comparison, and multi-class objective functions. In yet another embodiment, a Hidden Markov Model (HMM) may be used to account for the temporal order of frames, where the HMM may consist of several GMMs and the probabilities describing the transitions between them.

FIG. 4 provides an example of a tag arrangement comprising a plurality of extracted features along with a value distance/tolerance, where each feature value is expressed as a tolerable range for later comparison. Under one embodiment, each extracted audio feature is separately measured and collected as ranges (410A-420A) for template 400. Depending on the feature extracted, ranges may be combined, weighted, averaged and/or normalized for unit variance. Ranges are then set against value distances that are determined through any of Euclidean (e.g., 413A, 417A-419A), weighted Euclidean (e.g., 410A-412A, 414A), Kullback-Leibler distances (e.g., 415A, 416A) or others for tag creation/identification 425. In the example of FIG. 4, audio features relating to timbre 410 may include specific measurements directed to mean and variance of the spectral centroid, roll-off, flux, and or percentage of low/high energy frames. Timbre-related measurements may be taken across a plurality of audio signals to establish a set of ranges 410A for a particular tag (425). Subsequent measurements that are compared against these ranges (e.g., using weighted Euclidean distance), and which fall within the predetermined range, may be said to match a feature range for a particular tag. Additional features may include a first MFCC measurement 411, involving the mean and variance of a predetermined number of mel-frequency cepstral coefficients or number of dimensions (411A), and a concatenation of timbre and MFCC features 412, 412A.

Beat histogram features 413 may also be used to identify prominent beats, which may comprise amplitudes and periods of peaks in the histogram, a ratio between the peaks and

the sum of all peaks 413A. Pitch 414 may be derived from a histogram of pitches in an audio signal 414A, which may include periods and amplitudes of prominent peaks on a full semitone scale and/or octave independent scale. Additional MFCCs 415 may be estimated from short audio frames, where a Gaussian Mixture Model (GMM) may be trained to model them 415A. Loudness 416 may be measured from the sum of frequency bands distributed on a Bark scale, where a GMM may be trained on the loudness values 416A. Spectral histogram 417 may be formed from a derivative of raw tone features, where the number of loudness levels exceeding a predetermined threshold in each frequency may be counted 417A. A Periodicity histogram 418 may measure periodic beats 418A, or a fluctuation pattern 419 may be used to measure periodicities in a signal 419A. It is understood that the examples of FIG. 4 are merely illustrative, and that other features/techniques described herein may be used for creating tags 425 for template 400. Additionally, other techniques such as a multivariate autoregressive model 420 may be used to capture temporal correlations of MFCCs over relatively short (e.g., 1-2s) segments to produce feature vectors for each segment (420A). The vectors may be used individually or combined to compare for later comparison to new incoming audio features to identify audio features and characteristics.

Under a preferred embodiment, each of templates (400-400B) is comprised of one or more tags 425. In one embodiment, each tag is associated with a specific audio feature range (410A-420A). In another embodiment, each tag is associated with a plurality of audio feature ranges. For example, a tag relating to a genre, instrument or emotive descriptor may combine audio feature ranges from audio timbre (410A), beat (413A), loudness (416A) and spectral histogram (417A). In another example, the combined features may include audio timbre (410A), MFCC1 (411A), T+M (412A), and loudness (416A). In yet another example, combined features may include beat (413A) and periodicity histogram (418A). As can be appreciated by one skilled in the art, a large number of possibilities are available for developing and customizing tags for retrieval of semantic audio information/data.

Templates are preferably formed using a training process, where known audio signals are fed into a system such as the one illustrated in FIG. 2, and audio features are identified and tagged. For example, a collection of songs, known to be from a specific genre, have a certain number of audio features extracted, where audio feature ranges are determined for each template. The type and number of audio features used is not critical and may be left to the discretion of the designer. If more audio features are used, this will likely result in more accurate and/or granular semantic data. However, increasing the number of features increases the processing power needed to extract and tag audio features. As the features are extracted, they may be joined to form ranges for features, and/or normalized or concatenated to form one or more feature vectors that are subsequently tagged. Once the requisite features are tagged, they are formed as part of a template that is deemed representative of a specific genre (e.g., jazz, classical, rock, etc.). The same techniques may be used to form representative templates for instruments, emotive descriptors, etc. As each template is formed, it is stored in a database (e.g., SQL) for retrieval, processing and comparison to new incoming audio. These operations are preferably performed in a backoffice application (e.g., 108, 109) using Qt SQL libraries such as QSqlDatabase and QSqlQuery. The backoffice should also be usable with various engines, from a simple SQLite file to MySQL,



PostgreSQL, Oracle, Access DB files or any DB supporting ODBC (Open Data Base Connectivity protocol).

One of the advantages of this configuration is that individual song-vs.-song comparison is not required, which can greatly reduce the database storage need. Instead of requiring thousands or millions of audio data features, databases may contain just a few hundred templates that are representative of a specific audio characteristic. Such a configuration greatly simplifies the data needed for obtaining semantic information on audio. Of course, these audio features may be further decreased or increased, depending on the granularity needed.

Turning to FIG. 5, an exemplary comparison result is illustrated for an incoming audio signal that is processed and compared to an audio template described above. When a new audio signal is received, it may be processed according to a process described below in FIG. 6, and the resulting semantic audio signature is compared to a previously stored template created during a training process. As audio features are compared to templates, tagged audio features are identified and scored, and may further be aggregated into one or more score file histograms 500, where each file histogram 500 contains a score 501 relating to each respective feature. File 500 may consist of a single feature, or may contain a plurality of different features. In the example of FIG. 5, multiple features are contained in file 500, where features are related to various semantic information such as genre (classic jazz), instrumentation (acoustic drums, saxophone), style (swing), acoustical dynamics (dynamic, energetic) and emotive descriptors (happy). Again, the specific types and numbers of features are not critical and are left to the discretion of the designer. The resulting files are preferably time stamped and stored for later retrieval and processing.

FIG. 6 provides an example of new incoming audio (or test audio) received on a user device (e.g., 104-106), such as a cell phone, smart phone, personal computer, laptop, tablet, set-top-box, media box, and the like. The audio may be captured using a microphone that transduces the ambient audio into electrical form, or captured directly using a sound card, audio interface or the like. Incoming audio 601 is received and subjected to feature extraction 602 and feature integration 603, similar to the techniques described above in connection with FIG. 2. Depending on the device used, the number of extracted audio features may be increased or decreased, depending on the processing power and storage available. A semantic audio signature 604 is then formed from the extracted audio features, and stored on the user device. Under a preferred embodiment, the semantic audio signature is time stamped to indicate a time in which the signature was formed.

Semantic Signature 604 is then transmitted from the device via wired, wireless and/or cellular communication to a remote location, where the signature 604 is compared to audio templates 605, where tags are identified, scored and correlated. Under an alternate embodiment, the device may simply sample a time period of audio and transmit the sample via wired, wireless or cellular communication to a remote site for audio feature extraction, integration and semantic audio signature formation (604). Once tags are scored, they may be collected over a predetermined time period and processed for report generation. Unlike conventional audio signatures, semantic audio signatures may be taken over longer time intervals (e.g., 10-30 sec.), resulting in a saving of processing power.

FIGS. 7A and 7B illustrate a few examples of reports generated using the techniques described in FIG. 6. FIG. 7A illustrates a report 700 for a particular user ("User00001"),

where one semantic feature is monitored. In this case, the user's device is monitored to determine the type and/or genre of audio or music that the user was listening to at given times. Here, it can be determined that the user began his morning at 9:20 AM listening to talk programming (e.g., talk radio, podcast, etc.), and at 9:26 AM began listening to classical music. At 10:20 AM, the user listened to jazz, followed by classic rock at 11:20 AM, and returning back to talk programming at 12:00 PM. At 2:00 PM the user then listened to hard rock. FIG. 7B illustrates an example where multiple semantic features were used for the content discussed in FIG. 7A. Here, in addition to genre-related information (talk, classical, jazz, classic rock, hard rock), other semantic features, including instrumentation (woodwinds, saxophone, electric guitar), style (Baroque, conversational, cool, swing, confrontational, distortion), acoustical dynamics (aggressive, energetic) and emotive descriptors (happy, brooding) may be included as well.

It can be appreciated by those skilled in the art that the semantic information extracted from audio may provide additional and valuable information regarding user listening habits. Such information would be particularly valuable to those engaged in the audience measurement business to determine generic listening habits of users or panelists. Additionally, the semantic information may be used to establish "emotional profiles" for users and groups of users during the course of a day, week, month, year, etc. Demographic information may further be used to expand on these profiles to obtain demographically-related listening/emotional information.

Returning back to FIG. 2, semantic information may be used in conjunction with audio signatures that are processed independently via signature generation assembly 218. Signature generation assembly 218 may be based on analog and/or digital processing, and may generate audio signatures based on time and/or frequency components. The hardware components are generally known in the art and will not be discussed in detail for the sake of brevity. Exemplary configurations are disclosed in U.S. Pat. No. 5,436,653 titled "Method and System for Recognition of Broadcast Segments" issued Jul. 25, 1995 to Ellis et al., and U.S. Pat. No. 6,990,453, titled "System and Methods for Recognizing Sound and Music Signals in High Noise and Distortion" issued Jan. 24, 2006 to Li-Chun Wang et al., which are incorporated by reference in their entirety herein. Here, sampled audio segments stored in buffer 215 are forwarded to pre-processing 216, which may perform necessary windowing, gain control, filtering and transformation to the audio signal before feature extraction is performed. In block 217, audio signature feature extraction is performed in the frequency and/or time domain to generate audio signature 218.

FIG. 2A illustrates an exemplary technique in which an audio signature may be generated. In certain embodiments, when using data resulting from an FFT performed across a predetermined frequency range, the FFT data from an even number of frequency bands (for example, eight, ten, sixteen or thirty two frequency bands) spanning the predetermined frequency range are used two bands at a time during successive time intervals. FIG. 2A provides an example of how pairs of the bands are selected in these embodiments during successive time intervals where the total number of bands used is equal to ten. The selected bands are indicated by an "X".

When each band is selected, the energy values of the FFT bins within such band and such time interval are processed to form one bit of a signature. If there are ten FFT's for each



time interval of the audio signal, for example, the values of all bins of such band within the first five FFT's are summed to form a value "A" and the values of all bins of such band within the last five FFT's are summed to form a value "B". In the case of a received audio signal, the value A is formed from portions of the audio signal that were transmitted prior to those used to form the value B or which represent earlier portions of the audio signal relative to its time base.

To form a bit of the signature, the values A and B are compared. If B is greater than A, the bit is assigned a value "1" and if A is greater than or equal to B, the bit is assigned a value of "0". Thus, during each time interval, two bits of the signature are produced. Each bit of the signature is a representation of the energy content in the band represented thereby during a predetermined time period, and may be referred to as the "energy slope" thereof. Because any one energy slope is associated with a particular band, as opposed to being associated with a representation of energy content across a group of bands or between certain ones of various bands, the impact of fluctuations in the relative magnitudes of reproduced audio among frequency bands is virtually eliminated.

In certain embodiments, signatures may be extracted continuously. In such embodiments, information is obtained without a dependency on a triggering, predetermined event, or other type of prompting, and thus through uninterrupted information gathering, the signatures obtained will, necessarily, contain more information. For instance, this additional information is manifested in a signature, or portion thereof, that is formed of information as to how the audio signal changes over time as well as with frequency. This is in contrast to signature extraction occurring only upon prompting caused by a predetermined event and detection thereof, whereby information then obtained is only representative of the audio signal characterized within a certain isolated time frame.

Frequency bins or bands of different size may be employed to extract signatures and read codes. For example, relatively narrow bin sizes, such as 2, 4 or 6 Hz may be used to detect the presence of a component of an ancillary code, while signature extraction requires the use of wider bands, such as 30, 40 or 60 Hz to ensure that the band energy is sufficient to permit the extraction of a reliable signature or signature portion. Accordingly, in an advantageous embodiment of the invention that employs a time domain-to-frequency domain transformation that distributes the energy of an audio signal into a plurality of frequency bins or bands, the size or sizes of the bins or bands are each selected to have a first, relatively narrow frequency width. The energy values of such frequency bins or bands are processed to read an ancillary code therefrom. These energy values are also combined in groups of contiguous bins or bands (such as by addition) to produce frequency band values each representing an energy level within a frequency band comprising the respective group. Such frequency band values are then processed to extract a signature therefrom.

Another exemplary audio signature formation technique is described in A. Wang, "An industrial strength audio search algorithm," *Proceedings of the International Conference on Music Information Retrieval (ISMIR)*, pages 7-13, Baltimore, USA, 2003, where a spectrogram is formed from the audio, and spectrogram peaks on a time and frequency axis are used to generate a signature. Using spectrogram peaks may be advantageous since they are more likely to survive ambient noise, and further satisfy the property of linear superposition (i.e., a spectrogram peak analysis of music and noise together will contain spectral peaks due to the music

and the noise as if they were analyzed separately). Frequency peak points are identified and measured as a function of time from a predetermined time period (e.g., the beginning of a track). For pairs of peaks  $(t_1, f_1)$  and  $(t_2, f_2)$ , a fingerprint may be computed on a triplet of  $((t_2-t_1), f_1, (f_2-f_1))$ . Each number in the triplet may be quantized and the concatenated value is treated as the fingerprint. In one example, frequency data in a spectrogram may be divided into 256 levels linearly. Neighboring peaks may be considered in an adjacent frequency range of 64 units, and timing range of 64 units (when sampling rate of the audio signal is set to 8 KHz). The values  $((t_2-t_1), f_1, (f_2-f_1))$  thus may be represented as 6, 8 and 6 bits respectively to obtain a 20 bit fingerprint, where 20 fingerprints per second may be generated. Other configurations are possible for generating larger-sized fingerprints (e.g., 32-bit) as well.

Yet another exemplary audio signature formation technique is disclosed in U.S. patent application Ser. No. 13/228, 200, titled "Digital Signatures" to Samari et al., filed Sep. 8, 2011, which is incorporated by reference in its entirety herein. In this exemplary embodiment, audio signatures are generated on a time base (i.e., without transformation), where an envelope generator generates an envelope representation of a sampled audio data segment, and may (a) remove samples in the data segment which are of polarity opposite to the one polarity, or (b) to reverse the polarity of samples in the data segment which are of polarity opposite to the one polarity, to provide a polarized data segment. A sample value summer provides a portion sum value for successive portions of the polarized data segment comprising a predetermined plurality of samples. The portion sum value may represent the sum the values of the samples in the portion for the envelope representation. A threshold value generator then determines threshold values for successive blocks of the envelope representation, where each threshold value may be based on a par value of the portion sum values of the portions of the block concerned. Alternately, the threshold value generator determine a threshold value for each portion of the envelope representation, where each threshold value may be derived from (1) the threshold of the preceding portion, (2) the current portion sum value, and/or (3) a multiplication factor. An event detector may detect, as an event, a transition of a portion sum value across the threshold value for the block to which the portion concerned belongs, and a signature generator responds to the or a detected events to a generate a digital signature characteristic ("0" or "1", depending on whether or not a threshold is crossed) of the sampled data segment.

Formed audio signatures may subsequently be transmitted to central server(s) 109 via wired or wireless connection over a data network for audio signature matching as is known in the art. Once the signatures are matched, additional identification information (e.g., name of program, song, artist, performer, broadcaster, content provider, etc.) relating to audio 201 may be determined. The audio signatures may further be combined with semantic audio information to provide even more robust data.

By combining audio signatures with semantic information, this creates a powerful tool for assisting audience measurement entities in obtaining more information on the content that panelists are exposed to. Under the present disclosure, the semantic information can provide a deeper understanding of the underlying features of identified audio content. For example, a certain artist may perform songs spanning multiple genres. Using the techniques described herein, it can be automatically determined if certain genres by one artist are more popular than others. Similarly, it can



be automatically determined which ones of a plurality of artists of one genre are more popular than others. Furthermore, the techniques described herein may be used in television/streaming programming as well. For example, it may be determined that one or more panelists “tune out” a program when certain semantic features are present (e.g., confrontation, melancholy).

Turning to FIG. 8, a simplified example is provided where audio signatures are matched to identify content 802 (CONT1), and where the corresponding semantic information 801 is associated to it over a time base (starting from  $t=0$ ). Here, an 8-segment example is provided where content (CONT1) has a first kind of semantic information (SI1) present for time segments 1 and 2. At time segment 3, a second kind of semantic information (SI2) is present, followed by a third kind (SI3) for time segments 4 and 5. At time segments 6-8, the content returns to having semantic information (SI1). In a practical sense, the semantic information provides a “development” for the content over an extended time period. For example, if CONT1 is music, SI1 may represent a verse portion for a song, SI2 may represent a bridge, and SI3 may represent a chorus. Alternately, if CONT1 is a television program, SI1 may represent dialog, SI2 may represent the presence of dramatic music, and SI3 may represent a confrontational scene. It can be appreciated by those skilled in the art that many variations are possible under the present disclosure.

In addition to providing the development of content, semantic information may also be used to supplement audio signature data as well. Turning now to FIG. 9, a simplified example 900 is provided, where audio signatures are taken together with the semantic information on a time base to supplement audio signature detection. In this example, audio signature identification data 902 is associated with semantic information 903 according to timestamps 901 provided for each. While the timestamps themselves may provide an adequate basis for grouping audio signatures with semantic information, it is preferred that timestamp groupings are performed under a predetermined tolerance (+/-) to take into account possible time drift or skew that may occur during processing on a portable device. If the audio signatures and semantic audio signatures are not being processed simultaneously, a predetermined time delta may also be used to account for the time difference in which audio signatures and semantic audio signatures are generated.

In FIG. 9, a first audio signature (SIG0035) and related semantic audio information (Info1) is determined at time period 1. The audio signature (SIG0035) may be configured to provide specific information for the content (e.g., song), while the semantic information (Info1) may be configured to provide generalized information (e.g., genre, emotive descriptor). At time period 2, no audio signature was captured (“X”), which may occur because of excessive noise. However, semantic information (Info1) was determined for time period 2. Assuming that time periods 1 and 2 were sufficiently close in time, the presence of the same semantic information during those times would strongly suggest that that the same content (i.e., SIG0035) was being viewed. Accordingly, the content identification for SIG0035 from time period 1 may be extrapolated into time period 2.

Various embodiments disclosed herein provide devices, systems and methods for performing various functions using an audience measurement system that includes audio beaconing. Although specific embodiments are described herein, those skilled in the art recognize that other embodiments may be substituted for the specific embodiments shown to achieve the same purpose. As an example,

although terms like “portable” are used to describe different components, it is understood that other, fixed, devices may perform the same or equivalent functions. Also, while specific communication protocols are mentioned in this document, one skilled in the art would appreciate that other protocols may be used or substituted. This application covers any adaptations or variations of the present invention. Therefore, the present invention is limited only by the claims and all available equivalents.

The invention claimed is:

1. A processor-based method for producing supplemental information for audio signature data, the method comprising:

obtaining, by executing instructions with a processor, the audio signature data of a first time period, the audio signature data including data relating to at least one of time or frequency components representing a first characteristic of media;

obtaining, by executing instructions with the processor, first semantic audio signature data for the first time period, the first semantic audio signature data being a measure of generalized information representing characteristics of the media; and

storing, in a memory, the audio signature data of the first time period in association with a second time period when the processor determines, by executing instructions with the processor, that second semantic audio signature data for the second time period substantially matches the first semantic audio signature data for the first time period.

2. The method of claim 1, wherein the measure of generalized information representing characteristics of the media includes a timber feature having at least one of a mean of a spectral centroid, a variance of the spectral centroid, or a percentage of low/high energy frames.

3. The method of claim 1, wherein the measure of generalized information representing characteristics of the media includes a beat feature having at least one of an amplitude of peaks in a beat histogram, a period of peaks in the beat histogram, or a ratio between a peak and a sum of all peaks in the beat histogram.

4. The method of claim 1, wherein the measure of generalized information representing characteristics of the media includes a pitch feature having at least one of an amplitude of prominent peaks in a pitch histogram, or a period of peaks in the pitch histogram, where the pitch histogram is on a full semitone scale or an octave independent scale.

5. The method of claim 1, wherein the first semantic audio signature data includes at least one of a temporal feature, a spectral feature, a harmonic feature, or a rhythmic feature.

6. The method of claim 1, wherein the audio signature data is generated by transforming an audio signal of the media from a time domain to a frequency domain.

7. The method of claim 1, wherein the first semantic audio signature data is generated by transforming an audio signal of the media from a time domain to a frequency domain.

8. An apparatus for producing supplemental information for audio signature data, the apparatus including:

a processor to:

obtain the audio signature data of a first time period, the audio signature data including data relating to at least one of time or frequency components representing a first characteristic of media;

obtain first semantic audio signature data for the first time period, the first semantic audio signature data



## 21

being a measure of generalized information representing characteristics of the media; and  
 memory to store the audio signature data of the first time period in association with a second time period when the processor determines that second semantic audio signature data for the second time period substantially matches the first semantic audio signature data for the first time period.

9. The apparatus of claim 8, wherein the measure of generalized information representing characteristics of the media includes a timber feature having at least one of a mean of a spectral centroid, a variance of the spectral centroid, or a percentage of low/high energy frames.

10. The apparatus of claim 8, wherein the measure of generalized information representing characteristics of the media includes a beat feature having at least one of an amplitude of peaks in a beat histogram, a period of peaks in the beat histogram, or a ratio between a peak and a sum of all peaks in the beat histogram.

11. The apparatus of claim 8, wherein the measure of generalized information representing characteristics of the media includes a pitch feature having at least one of an amplitude of prominent peaks in a pitch histogram, or a period of peaks in the pitch histogram, where the pitch histogram is on a full semitone scale or an octave independent scale.

12. The apparatus of claim 8, wherein the first semantic audio signature data includes at least one of a temporal feature, a spectral feature, a harmonic feature, or a rhythmic feature.

13. The apparatus of claim 8, wherein the audio signature data is generated by transforming an audio signal of the media from a time domain to a frequency domain.

14. The apparatus of claim 8, wherein the first semantic audio signature data is generated by transforming an audio signal of the media from a time domain to a frequency domain.

15. A processor-based method for producing supplemental information for audio signature data, the method comprising:

## 22

obtaining, by executing instructions with a processor, the audio signature data at an input from a data network, the audio signature data received from a device, the audio signature data including data relating to at least one of time or frequency components representing a first characteristic of media;

obtaining, by executing instructions with the processor, semantic audio signature data at the input from the data network, the semantic audio signature data received from the device, the semantic audio signature data being a measure of at least one of a timber feature, a beat feature or a pitch feature relating to a second characteristic of the media;

associating, by executing instructions with the processor, the semantic audio signature data to the audio signature data; and

processing, by executing instructions with the processor, the associated semantic audio signature data and audio signature data to determine a change in the second characteristic relative to the first characteristic.

16. The method of claim 15, wherein the timber feature includes one of a mean of a spectral centroid, a variance of the spectral centroid, or a percentage of low/high energy frames.

17. The method of claim 15, wherein the beat feature includes one of an amplitude of peaks in a beat histogram, a period of peaks in the beat histogram, or a ratio between a peak and a sum of all peaks in the beat histogram.

18. The method of claim 15, wherein the pitch feature includes one of an amplitude of prominent peaks in a pitch histogram, or a period of peaks in the pitch histogram, where the pitch histogram is on a full semitone scale or an octave independent scale.

19. The method of claim 15, wherein the semantic audio signature data includes at least one of a temporal feature, a spectral feature, a harmonic feature, or a rhythmic feature.

20. The method of claim 15, wherein the semantic audio signature data is generated by transforming an audio signal of the media from a time domain to a frequency domain.

\* \* \* \* \*